



Performance of Future High-End Computers

David H. Bailey

Chief Technologist, Computational Research Dept

Lawrence Berkeley National Laboratory

<http://crd.lbl.gov/~dhbailey>



Laplace Anticipates Modern High-End Computers



“An intelligence knowing all the forces acting in nature at a given instant, as well as the momentary positions of all things in the universe, would be able to comprehend in one single formula the motions of the largest bodies as well as of the lightest atoms in the world, provided that its intellect were sufficiently powerful to subject all data to analysis; to it nothing would be uncertain, the future as well as the past would be present to its eyes.”

-- Pierre Simon Laplace, 1773



Who Needs High-End Computers?



Expert predictions:

- ◆ (c. 1945) Thomas J. Watson (CEO of IBM):
“World market for maybe five computers.”
- ◆ (c. 1975) Seymour Cray:
“Only about 100 potential customers for Cray-1.”
- ◆ (c. 1977) Ken Olson (CEO of DEC):
“No reason for anyone to have a computer at home.”
- ◆ (c. 1980) IBM study:
“Only about 50 Cray-1 class computers will be sold per year.”

Present reality:

- ◆ Many homes now have 5 Cray-1 class computers.
- ◆ Latest PCs outperform 1988-era Cray-2.



Evolution of High-End Computing Technology



1950	Univac-1	1 Kflop/s (10^3 flop/sec)
1965	IBM 7090	100 Kflop/s (10^5 flop/sec)
1970	CDC 7600	10 Mflop/s (10^7 flop/sec)
1976	Cray-1	100 Mflop/s (10^8 flop/sec)
1982	Cray X-MP	1 Gflop/s (10^9 flop/sec)
1990	TMC CM-2	10 Gflop/s (10^{10} flop/sec)
1995	Cray T3E	100 Gflop/s (10^{11} flop/sec)
2000	IBM SP	1 Tflop/s (10^{12} flop/sec)
2002	Earth Simulator	40 Tflop/s (4×10^{12} flop/sec)

We are on-track to achieve 1 Pflop/s before 2010.



Life Cycle of Scientific Applications

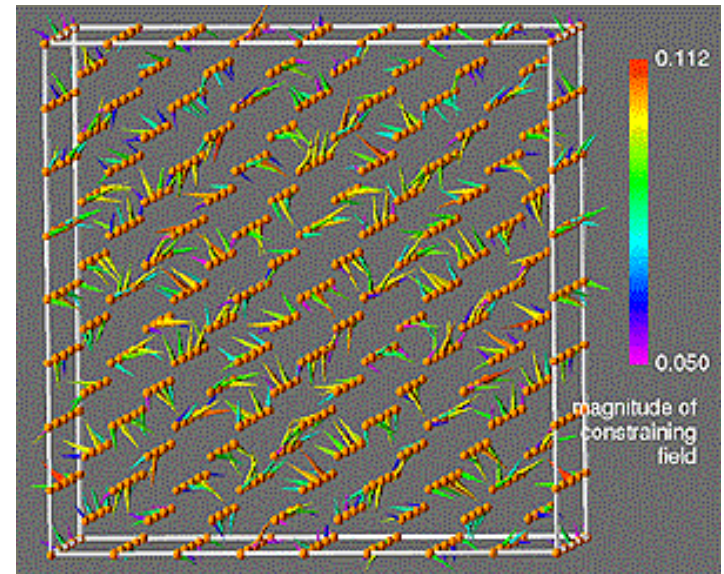
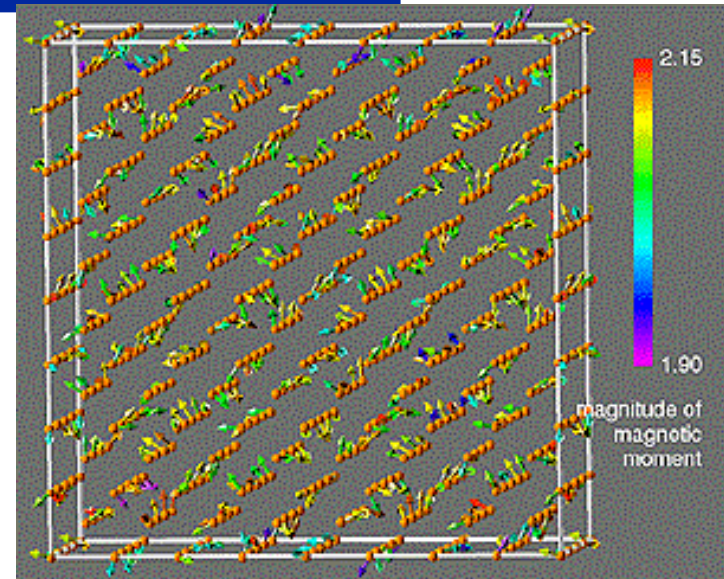


- ◆ Infeasible – much too expensive to consider.
- ◆ First sketch of possible computation.
- ◆ First demo on state-of-the-art highly parallel system.
- ◆ Code is adapted for production large-scale runs.
- ◆ Code runs on a shared memory multiprocessor.
- ◆ Code runs on a single-CPU workstation.
- ◆ Code runs on personal computer system.
- ◆ Code is embedded in web-based facility.
- ◆ Code is embedded in hand-held application.

- ◆ 6000-CPU IBM SP: 10 Tflop/s (10 trillion flops/sec).
- ◆ Currently the world's 5th most powerful computer.



- 1024-atom first-principles simulation of metallic magnetism in iron was 1998 Gordon Bell Prize winner -- first real scientific simulation to top 1Tflop/s.
- 2016-atom simulation now runs on the NERSC-3 system at 2.46 Tflop/s.





Materials Science Requirements



Electronic structures:

- ◆ Current: ~300 atom: 0.5 Tflop/s, 100 Gbyte memory.
- ◆ Future: ~3000 atom: 50 Tflop/s, 2 Tbyte memory.

Magnetic materials:

- ◆ Current: ~2000 atom: 2.64 Tflop/s, 512 Gbytes memory.
- ◆ Future: hard drive simulation: 30 Tflop/s, 2 Tbyte memory.

Molecular dynamics:

- ◆ Current: 10^9 atoms, ns time scale: 1 Tflop/s, 50 Gbyte mem.
- ◆ Future: alloys, us time scale: 20 Tflop/s, 4 Tbyte memory.

Continuum solutions:

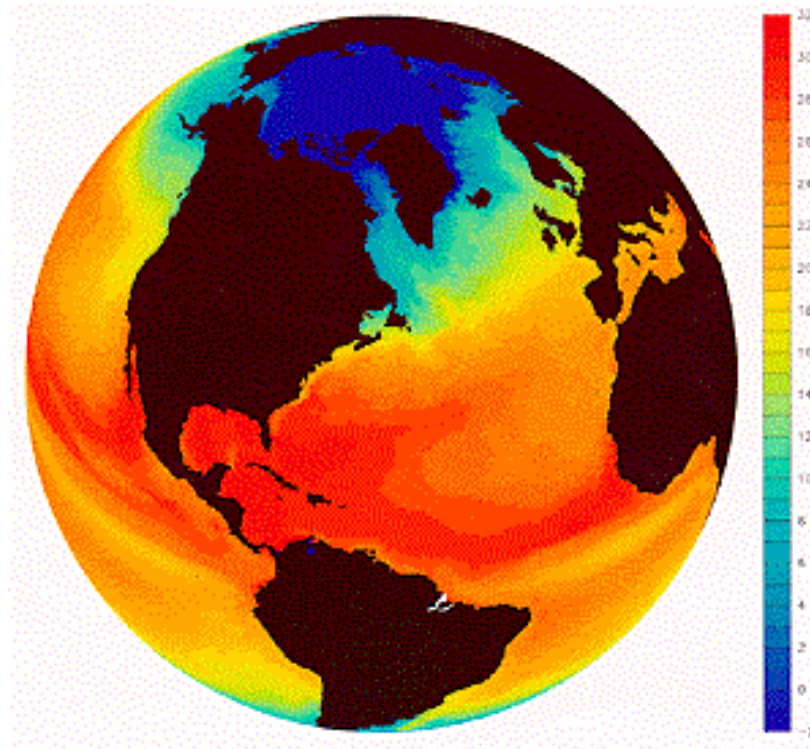
- ◆ Current: single-scale simulation: 30 million finite elements.
- ◆ Future: multiscale simulations: 10 x current requirements.



DOE Applications: Environmental Science



Parallel climate model (PCM) simulates long-term global warming.





Climate Modeling Requirements



Current state-of-the-art:

- ◆ Atmosphere: 1 x 1.25 deg spacing, with 29 vertical layers.
- ◆ Ocean: 0.25 x 0.25 degree spacing, 60 vertical layers.
- ◆ Currently requires 52 seconds CPU time per simulated day.

Future requirements (to resolve ocean mesoscale eddies):

- ◆ Atmosphere: 0.5 x 0.5 deg spacing.
- ◆ Ocean: 0.125 x 0.125 deg spacing.
- ◆ Computational requirement: 17 Tflop/s.

Future goal: resolve tropical cumulus clouds:

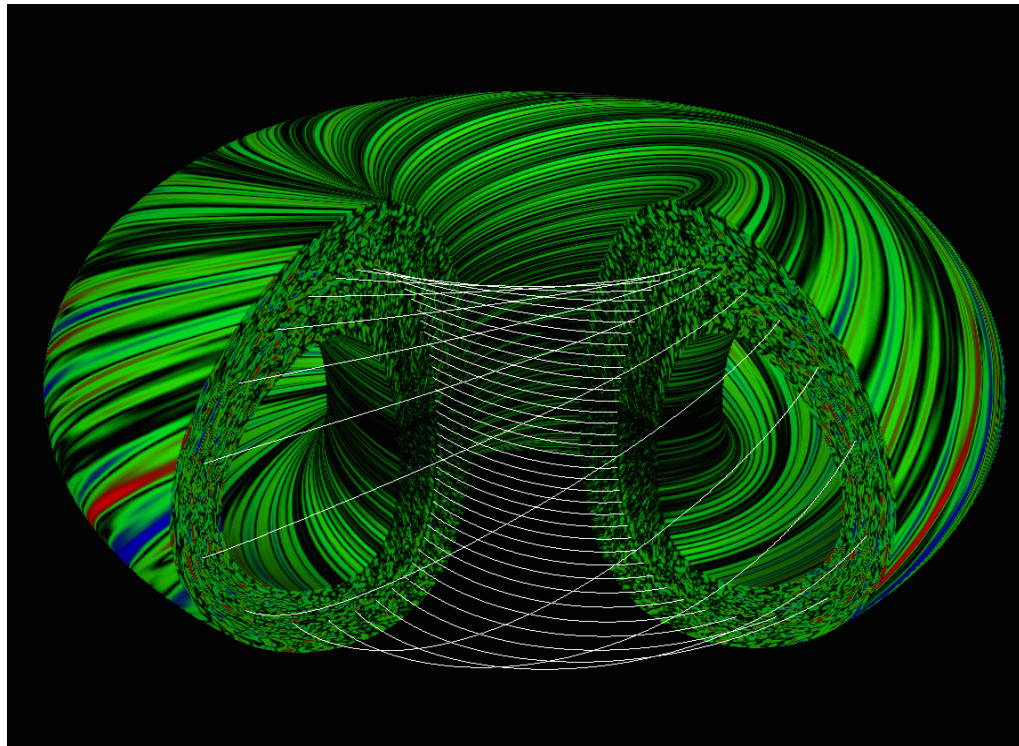
- ◆ 2 to 3 orders of magnitude more than above.



DOE Applications: Fusion Energy



Computational simulations help scientists understand turbulent plasmas in nuclear fusion reactor designs.



Tokamak simulation -- ion temperature gradient turbulence in ignition experiment:

- ◆ Grid size: $3000 \times 1000 \times 64$, or about 2×10^8 gridpoints.
- ◆ Each grid cell contains 8 particles, for total of 1.6×10^9 .
- ◆ 50,000 time steps required.
- ◆ Total cost: 3.2×10^{17} flop/s, 1.6 Tbyte.

All-Orders Spectral Algorithm (AORSA) – to address effects of RF electromagnetic waves in plasmas.

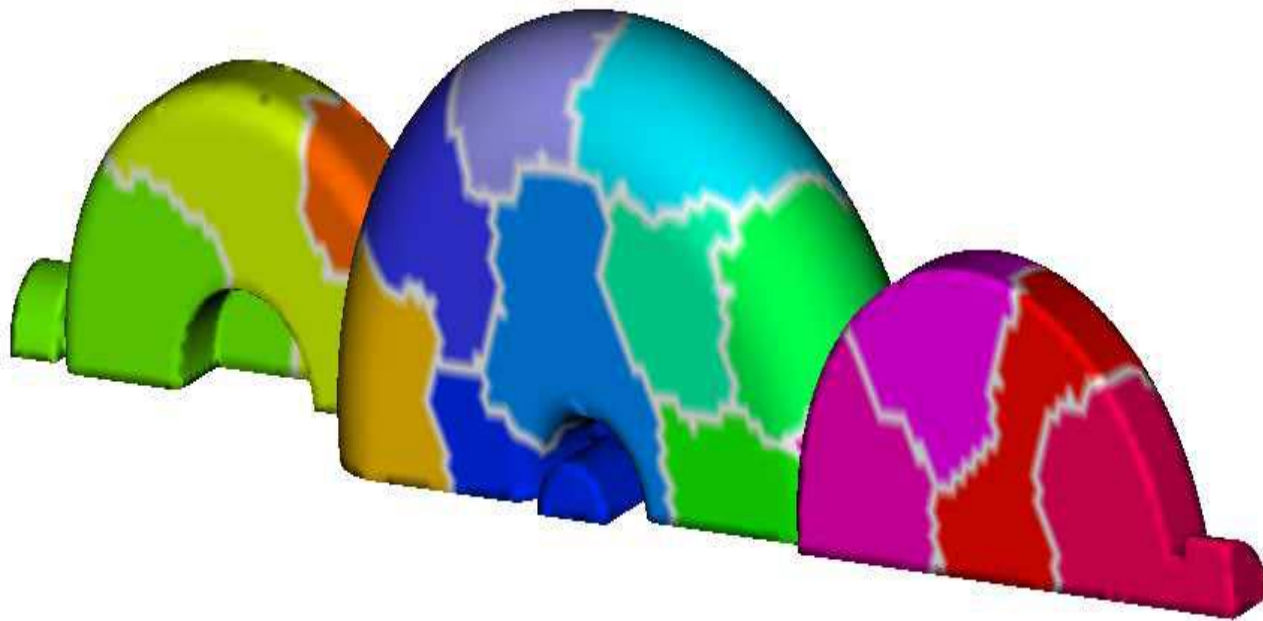
- ◆ $120,000 \times 120,000$ complex linear system.
- ◆ 230 Gbyte memory.
- ◆ 1.3 hours on 1 Tflop/s.
- ◆ $300,000 \times 300,000$ linear system requires 8 hours.
- ◆ Future: $6,000,000 \times 6,000,000$ system (576 Tbyte memory), 160 hours on 1 Pflop/s system.



NERSC/DOE Applications: Accelerator Physics



Simulations are being used to design future high-energy physics research facilities.





Accelerator Modeling Requirements



Current computations:

- ◆ 1283 to 5123 cells, or 40 million to 2 billion particles.
- ◆ Currently requires 10 hours on 256 CPUs.

Future computations:

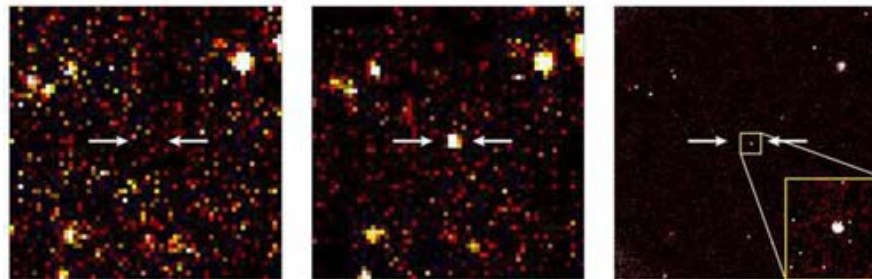
- ◆ Modeling intense beams in rings will be 100 to 1000 times more challenging.



DOE Applications: Astrophysics and Cosmology



- ◆ The oldest, most distant Type 1a supernova confirmed by computer analysis at NERSC.
- ◆ Supernova results point to an accelerating universe.
- ◆ Analysis at NERSC of cosmic microwave background data shapes concludes that geometry of the universe is flat.





Astrophysics Requirements



Supernova simulation:

- ◆ Critical need to better understand Type 1a supernovas, since these are used as “standard candles” in calculating distances to remote galaxies.
- ◆ Current models are only 2-D.
- ◆ Initial 3-D model calculations will require 2,000,000 CPU-hours per year, on jobs exceeding 256 Gbyte memory.
- ◆ Future calculations 10 to 100 times as expensive.

Analysis of cosmic microwave background data:

- | | | |
|----------------------|----------------------------|---------------|
| ◆ MAXIMA data | 5.3×10^{16} flops | 100 Gbyte mem |
| ◆ BOOMERANG data | 1.0×10^{19} flops | 3.2 Tbyte mem |
| ◆ Future MAP data | 1.0×10^{20} flops | 16 Tbyte mem |
| ◆ Future PLANCK data | 1.0×10^{23} flops | 1.6 Pbyte mem |

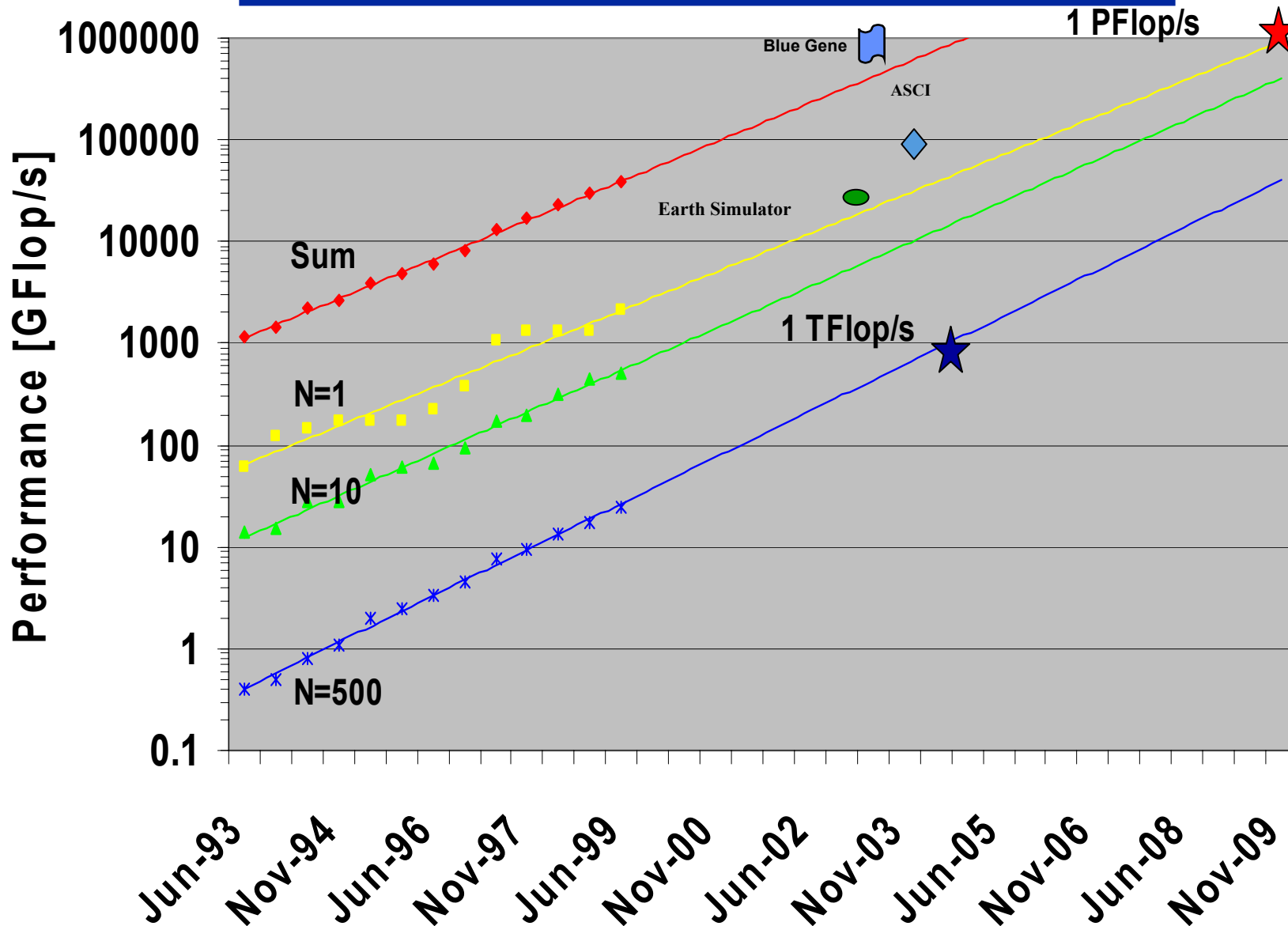


Applications for Petaflops Computers



- ◆ Weather forecasting.
- ◆ Business data mining.
- ◆ DNA sequence analysis.
- ◆ Protein folding simulations.
- ◆ Inter-species DNA analyses.
- ◆ Medical imaging and analysis.
- ◆ Nuclear weapons stewardship.
- ◆ Multiuser immersive virtual reality.
- ◆ National-scale economic modeling.
- ◆ Climate and environmental modeling.
- ◆ Molecular nanotechnology design tools.
- ◆ Cryptography and digital signal processing.

Top500 Trends





The Japanese Earth Simulator System



System design:

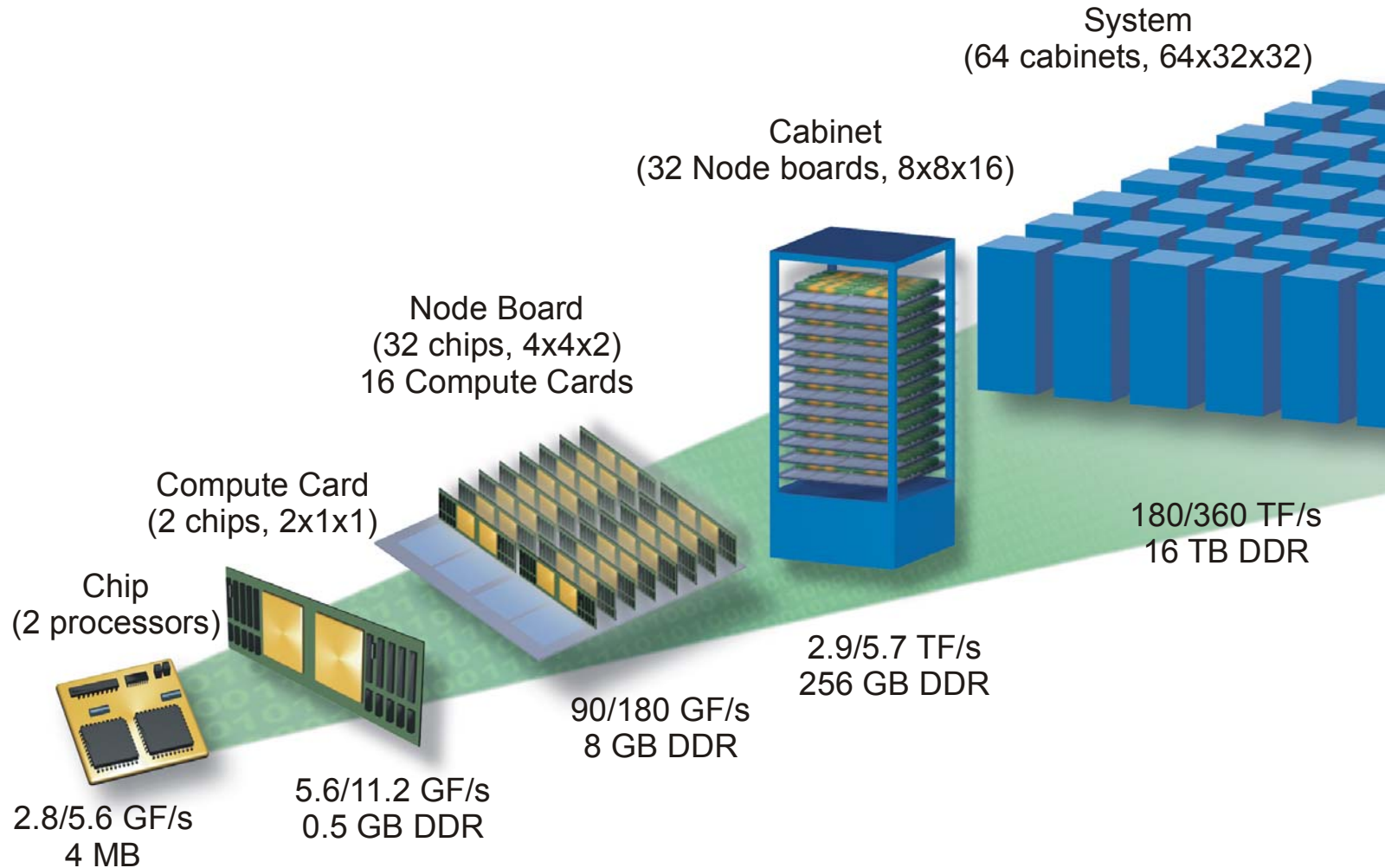
- ◆ Architecture: Crossbar-connected multi-proc vector system.
- ◆ Performance: 640 nodes x 8 proc per node x 8 Gflop/s per proc = 40.96 Tflop/s peak
- ◆ Memory: 640 nodes x 16 Gbyte per node = 10.24 Tbyte.

Sustained performance:

- ◆ Global atmospheric simulation: 26.6 Tflop/s.
- ◆ Fusion simulation (all HPF code): 12.5 Tflop/s.
- ◆ Turbulence simulation (global FFTs): 12.4 Tflop/s.



IBM's Blue Gene/L System





Other Future High-End Designs



◆ Processor in memory

- ◆ Currently being pursued by a team headed by Prof. Thomas Sterling of Cal Tech.
- ◆ Seeks to design a high-end scientific system based on special processors with embedded memory.
- ◆ Advantage: significantly greater processor-memory bandwidth.

◆ Streaming supercomputer

- ◆ Currently being pursued by a team headed by Prof. William Dally of Stanford.
- ◆ Seeks to adapt streaming processing technology, now used in game market, to scientific computing.
- ◆ Projects 200 Tflop/s, 200 Tbyte system will cost \$10M in 2007.

Gordon Moore, 1965:

"The complexity for minimum component costs has increased at a rate of roughly a factor of two per year... Certainly over the short term this rate can be expected to continue, if not to increase. Over the longer term, the rate of increase is a bit more uncertain, although there is no reason to believe it will not remain nearly constant for at least 10 years."
[Electronics, Apr. 19, 1965, pg. 114-117.]

Gordon Moore, 2003:

No end in sight – Moore's Law will continue for at least another ten years.

Moore's Law may even **accelerate** in the future.



Beyond Silicon: Sooner Than You Think



Nanotubes:

- ◆ Can function as conductors, memory and logic devices.
- ◆ Nantero, a venture-funded firm, has devices in development.

Molecular self-assembly:

- ◆ Researchers at HP have created a memory device with crisscrossing wires 2 nm wide, 9 nm apart.
- ◆ 1994 goal: 1000 bits/ μ^2 (compared to 10 bits/ μ^2 for DRAM).
- ◆ Zettacore, a venture-funded firm, is pursuing related ideas.

Molecular electronics:

- ◆ Researchers Mitre and UCLA have demonstrated organic molecules that act as electronic logic devices.

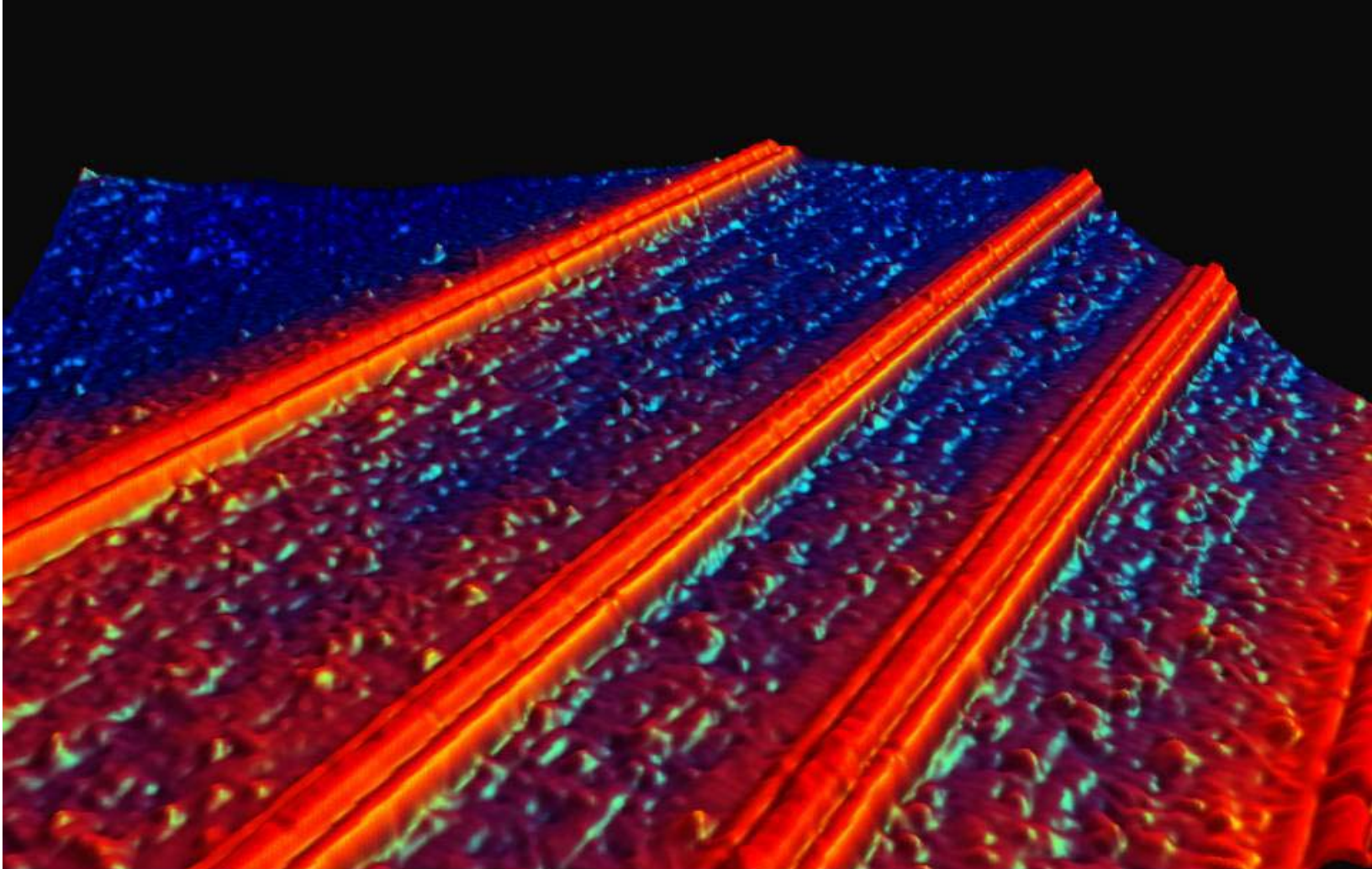
Atomic force microscopes (AFMs):

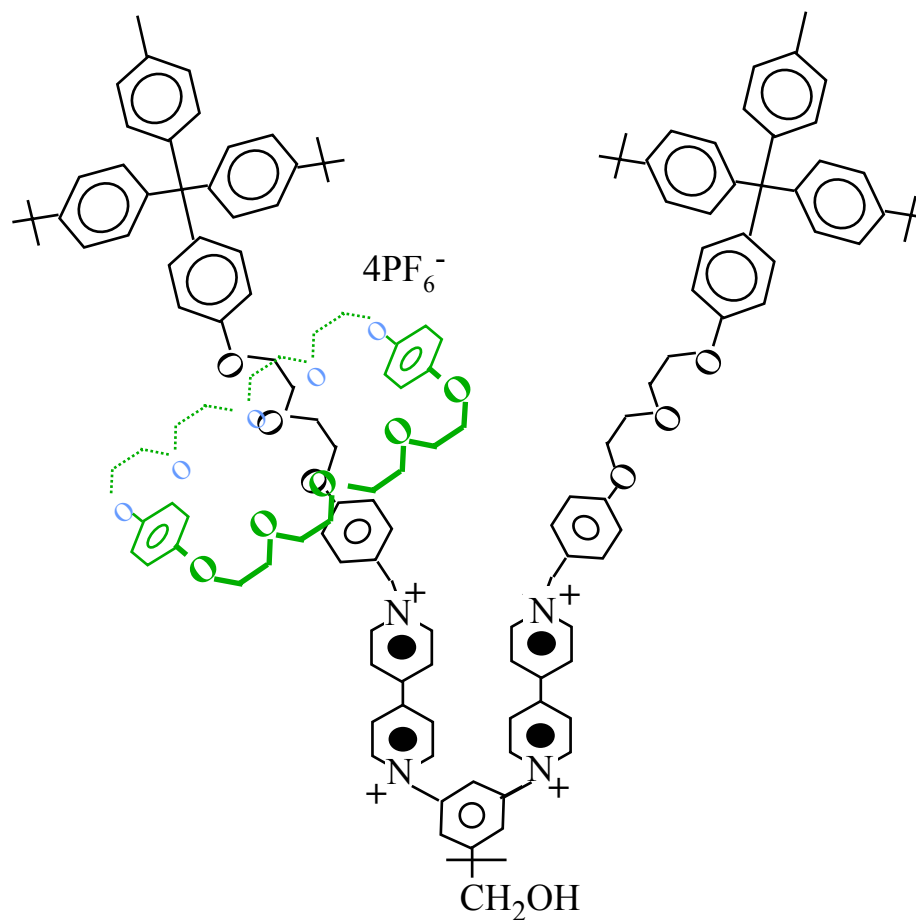
- ◆ Arrays of AFM tips can etch patterns 5nm lines on silicon.



Self-Assembled Wires 2nm Wide

[P. Kuekes, S. Williams, HP Labs]







Fundamental Device Limits



Assume:

- ◆ Power dissipation of 1 watt at room temperature,
- ◆ A spatial volume of 1 cm³.

Q: How many bit operation/second can be performed by a nonreversible computer executing Boolean logic?

A: $P/kT \log(2) = 3.5 \times 10^{20}$ bit ops/s

Q: How many bits/second can be transferred?

A: $\sqrt{cP/kTd} = 10^{18}$ bit/s

“There’s plenty of room at the bottom” -- Richard Feynman, 1959.



Research Questions for Future High-End Computing



- ◆ Can systems be designed with 10,000 to 100,000+ CPUs, with acceptable performance scaling?
- ◆ Will we be able to expose 10^8 -way concurrency in every significant step of major computations?
- ◆ Will 128-bit floating-point arithmetic be required?
- ◆ Will existing system software scale to this level?
- ◆ Will new programming models and/or languages be required?



The Performance Evaluation Research Center (PERC)



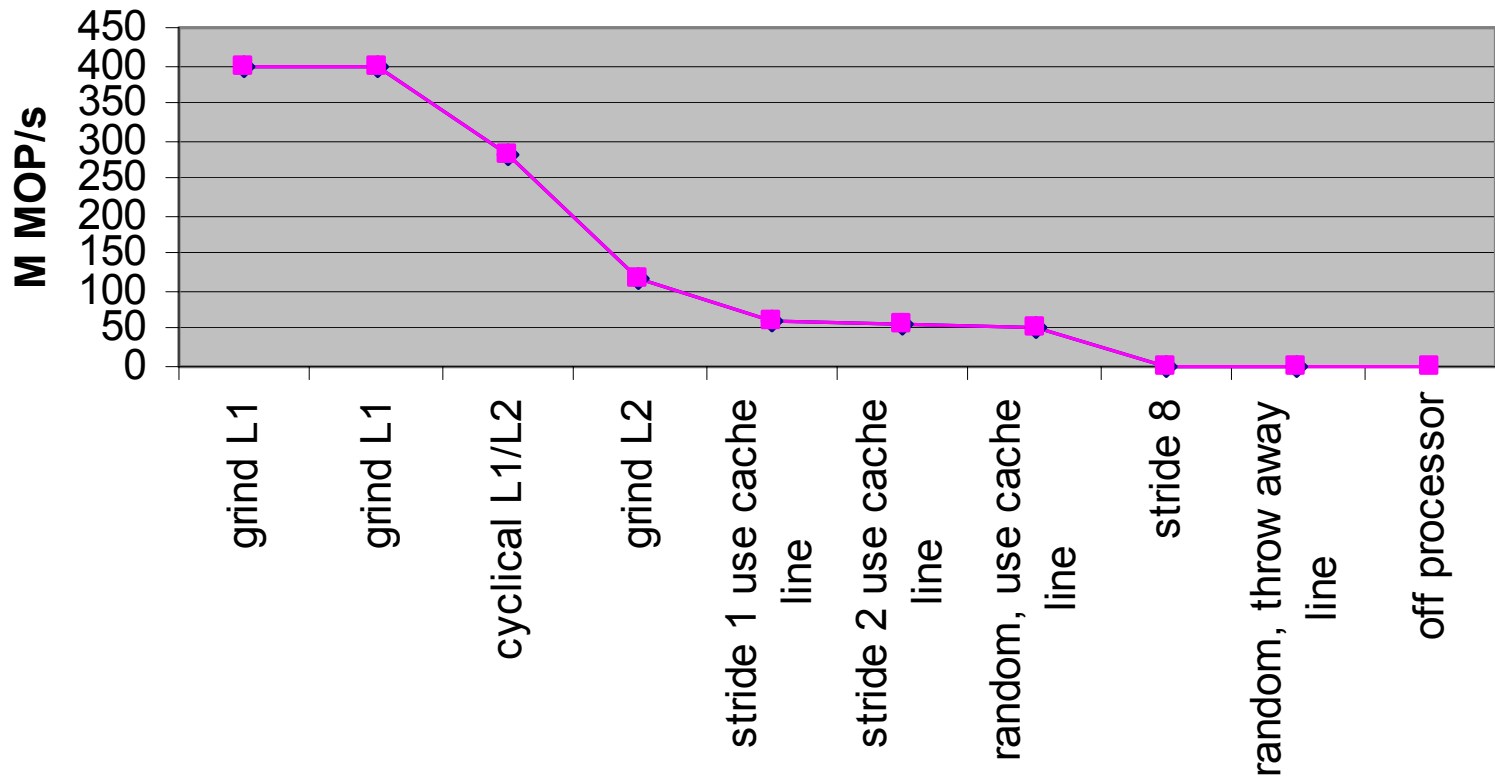
One of five Integrated Software Infrastructure Centers funded through the DoE SciDAC program.

Research thrusts:

- ◆ Develop improved tools for performance monitoring and code tuning.
- ◆ Study the performance characteristics of specific large-scale scientific codes.
- ◆ Develop tools and techniques for performance modeling.
- ◆ Develop semi-automatic facilities for improving performance.

Memory Access Patterns (MAPS)

MAPS Alpha 21264





User Tools: SvPablo



svPablo

Project Instrument View GenCallGraph Help

Project Description: PCTM on IBM-SP (seaborg)

Source Files:
fod.F
fod_setup.F

Performance Contexts:
IBM-SP, 64 procs, other Metrics
IBM-SP, 64 procs, other Metrics, 10 days
IBM-SP, 32 procs, other Metrics, 10 days
IBM-SP, 16 procs, other Metrics, 10 days
IBM-SP, 8 procs, other Metrics, 10 days
IBM-SP, 4 procs, other Metrics, 10 days

Routines in Source File
fod
fod_setup
fod_timer_clear
fod_timer_start
fod_init

Routines in Performance Data
com3
oceanstep
icestep
mpi_c
mpi_c

Source File: /u0/cmendes/PCTM/pam2/src/sources/fod.F

Specific Metric
Call Statistics count:
240.0000 -- com3
Dismiss Help

Specific Metric
Call Statistics Duration:
211.2399 -- com3
Dismiss Help

Specific Metric
HW Statistics by Line Floating Point Instructions:
12295122047.0000 -- com3
Dismiss Help

Specific Metric
HW Statistics by Line Load Misses in D1:
300348320.0000 -- com3
Dismiss Help

Specific Metric
HW Statistics by Line Branch Instructions:
6944481284.0000 -- com3
Dismiss Help

Specific Metric
HW Statistics by Line TLB misses:
151118598.0000 -- com3
Dismiss Help

Legend: Source Code Metrics

Column 1: Call Statistics count
240
10

Column 2: Call Statistics Duration
211.24
119.525

Column 3: Loop Statistics count
0
0

Column 4: Loop Statistics Duration
0
0

Column 5: HW Statistics by Line Floating Point Instr
1.65722e+10
0

Column 6: HW Statistics by Line Load Misses in D1
8.62196e+08
2.46132e+08

Column 7: HW Statistics by Line Branch Instructions
6.94448e+09
0

Column 8: HW Statistics by Line Load Instructions
3.10653e+10
0

Column 9: HW Statistics by Line Instruction Cache M
9.55654e+07
6.80955e+07

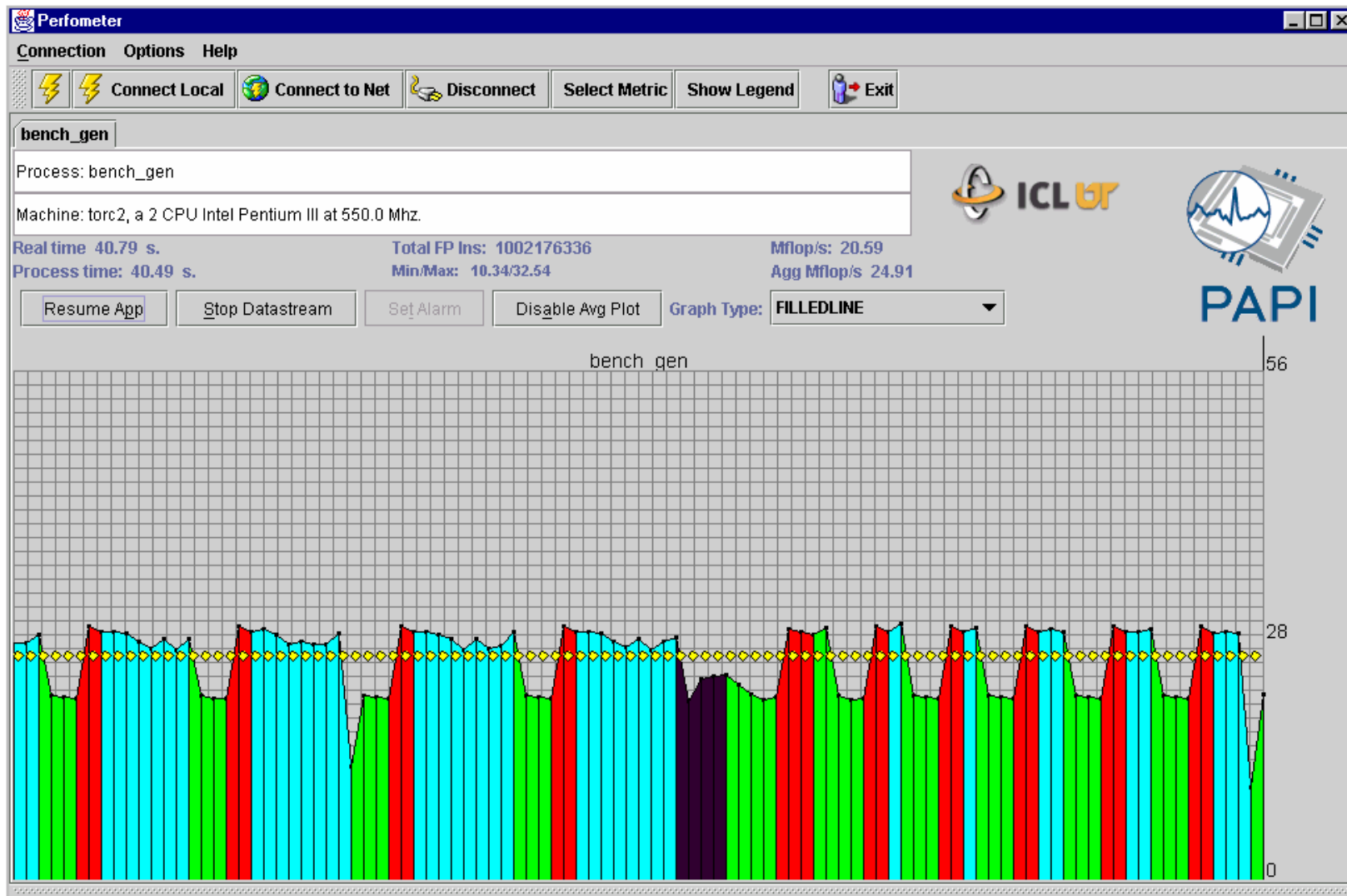
Column 10: HW Statistics by Line TLB misses
1.51119e+08
1.00566e+07

Dismiss Help

Instrument/Clear Line View Line Data



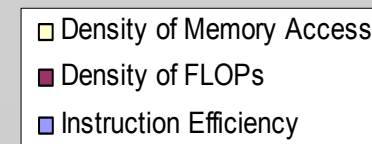
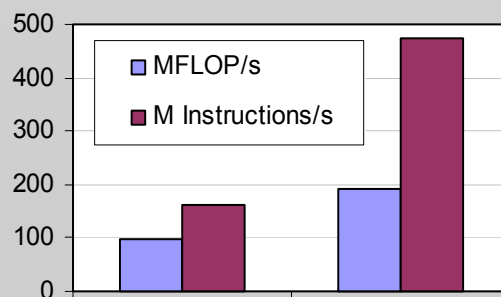
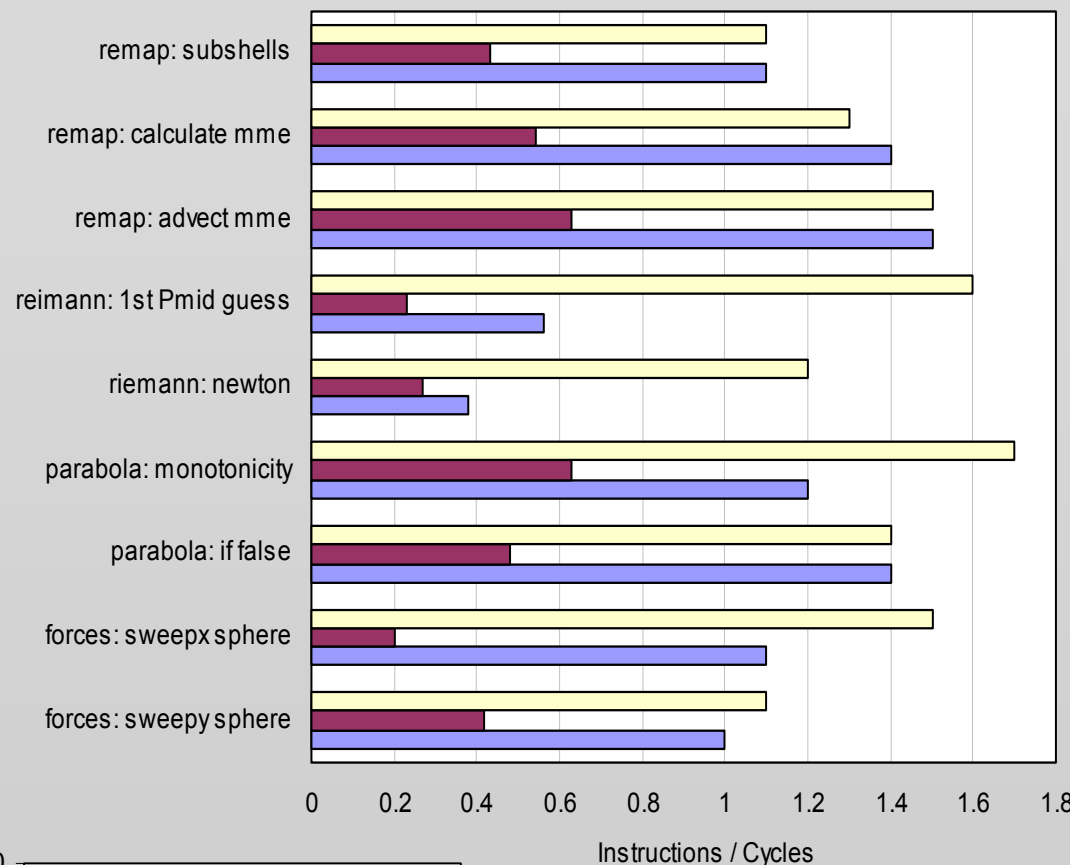
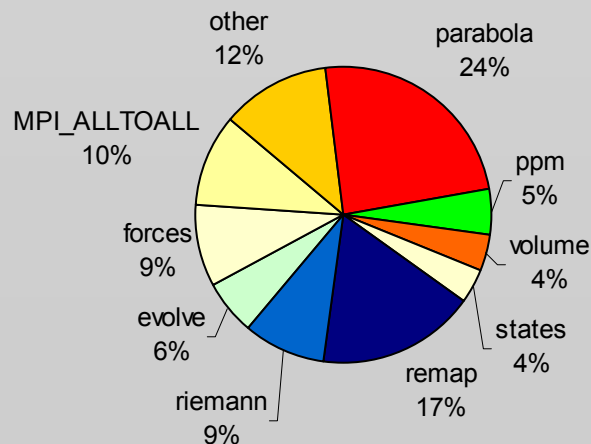
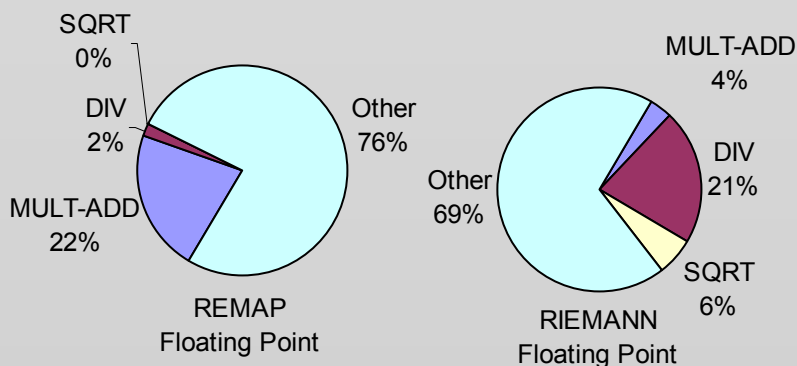
PAPI Perfometer Interface



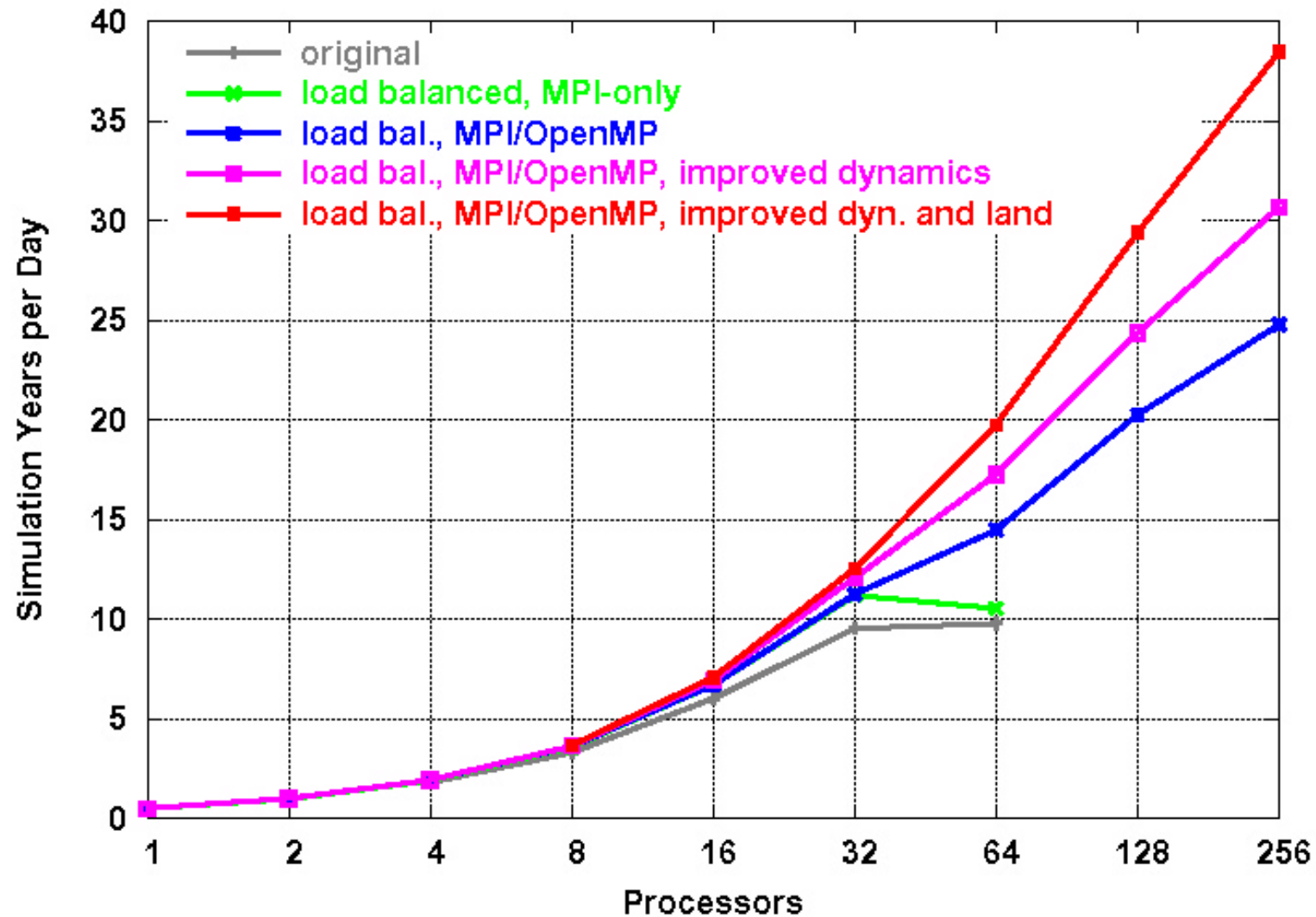
Performance Analysis

EVH1 (high-energy physics)

Aggregate performance measures over all tasks for a .1 simulation-second run. Collected with PAPI on an IBM SP (Nighthawk II / 375MHz).



Improvement to Climate Modeling Code





PERC Performance Modeling



- ◆ *Application signature* tools characterize applications independent of the machine where they execute.
- ◆ *Machine signature* tools characterize computer systems, independent of the applications.
- ◆ *Convolution* tools combine application and machine signatures to provide accurate performance models.
- ◆ *Statistical models* find approximate performance models based on easily measured performance data.

# CPUs	Real Time	Predicted Time	% Error
2	31.78	31.82	0.13
4	29.07	31.27	7.57
8	36.13	33.72	6.67
64	44.91	43.91	2.23
96	48.87	47.15	3.52
128	52.88	52.46	0.79



Future Challenges in the Performance Research Field



- ◆ Scaling performance monitoring tools to many thousands of processors.
- ◆ Handling the exploding volume of trace data.
- ◆ Visualizing performance results.
- ◆ Understanding the behavior, via simulation and modeling, of interprocessor networks.
- ◆ Developing performance modeling tools accurate enough to predict performance in the 10,000-100,000 processor arena.
- ◆ Extension of tools and modeling techniques to cover emerging architectures (Cray X1, ESS, PIM, BG/L).